# Fast-growing datasets meet slowly improving memory bandwidth and latency

## Growth of GenBank
### (1982 - 2005)



Doubling time for sequence databases is currently ~18 months

According to Moore's Law, doubling time for processor speed is ~18 months.

Time for doubling of bandwidth to memory and to disk = 2.7 years*
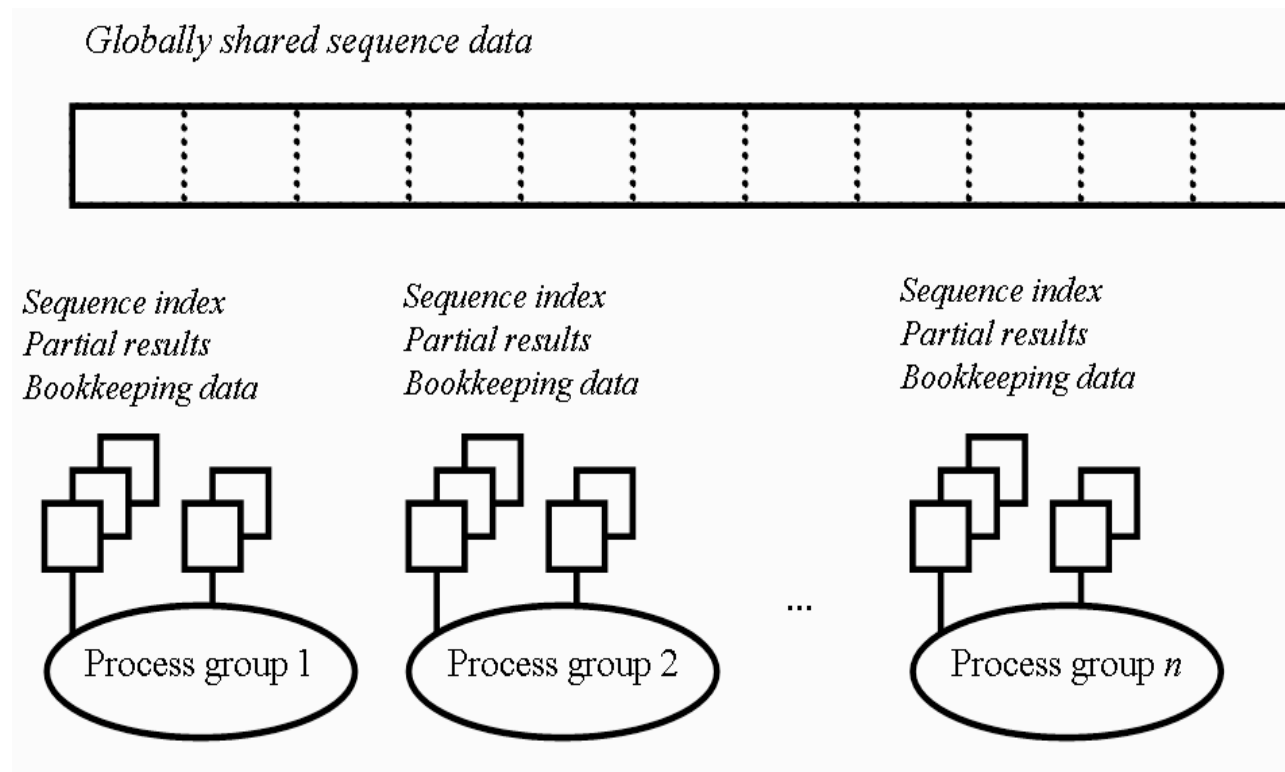
During this same time, memory *latency* only improves by 20%, and disk *latency* only improves by 30%*

*source: Patterson DA, "Latency Lags Bandwidth: Recognizing the chronic imbalance between bandwidth and latency, and how to cope with it", *Comm. ACM*. 47(10): **2004**, 71-75
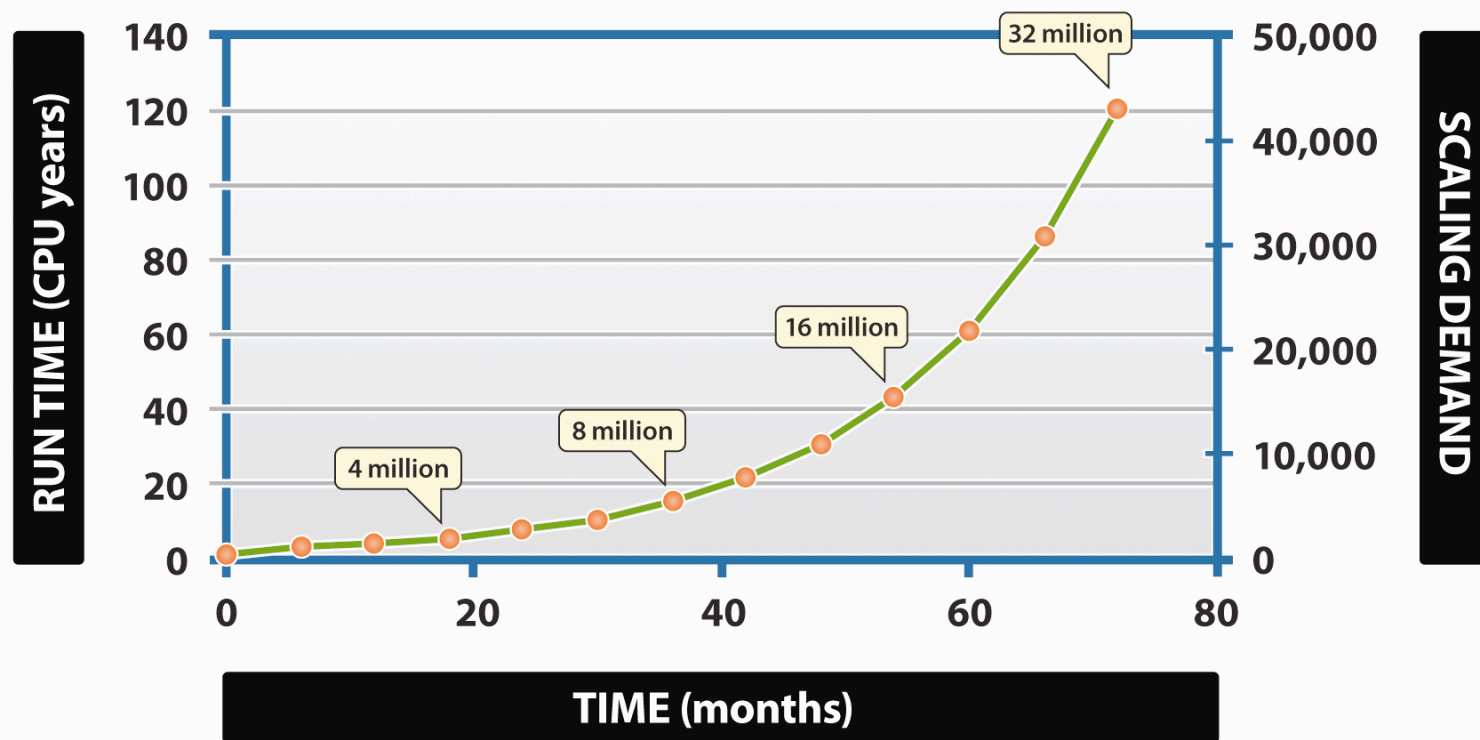
# ScalaBLAST
# Key: Memory Management

- Use large aggregate memory to share a single copy of the target database
- Hide latency by prefetching sequences in blocks.
- Each process group operates on independent query sets



Globally shared sequence data

Sequence index
Partial results
Bookkeeping data

Sequence index
Partial results
Bookkeeping data

Sequence index
Partial results
Bookkeeping data

Process group 1
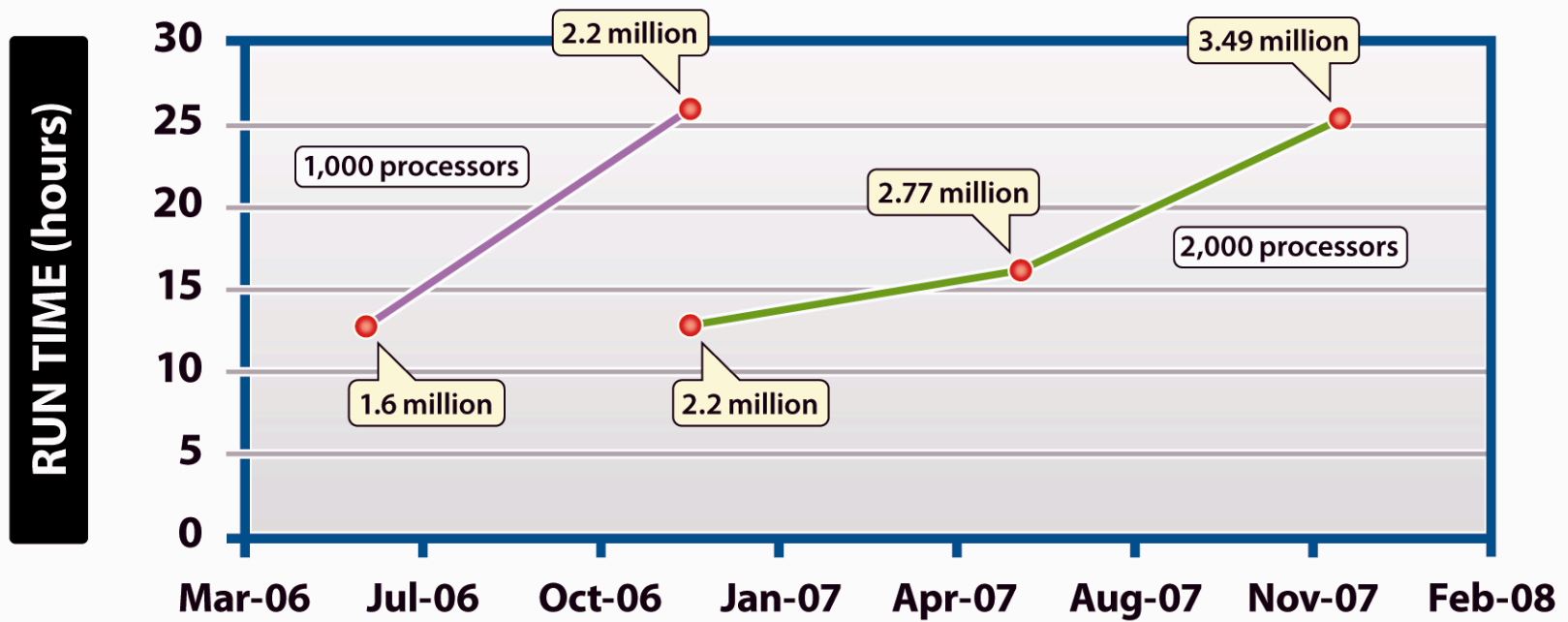
Process group 2

...

Process group *n*
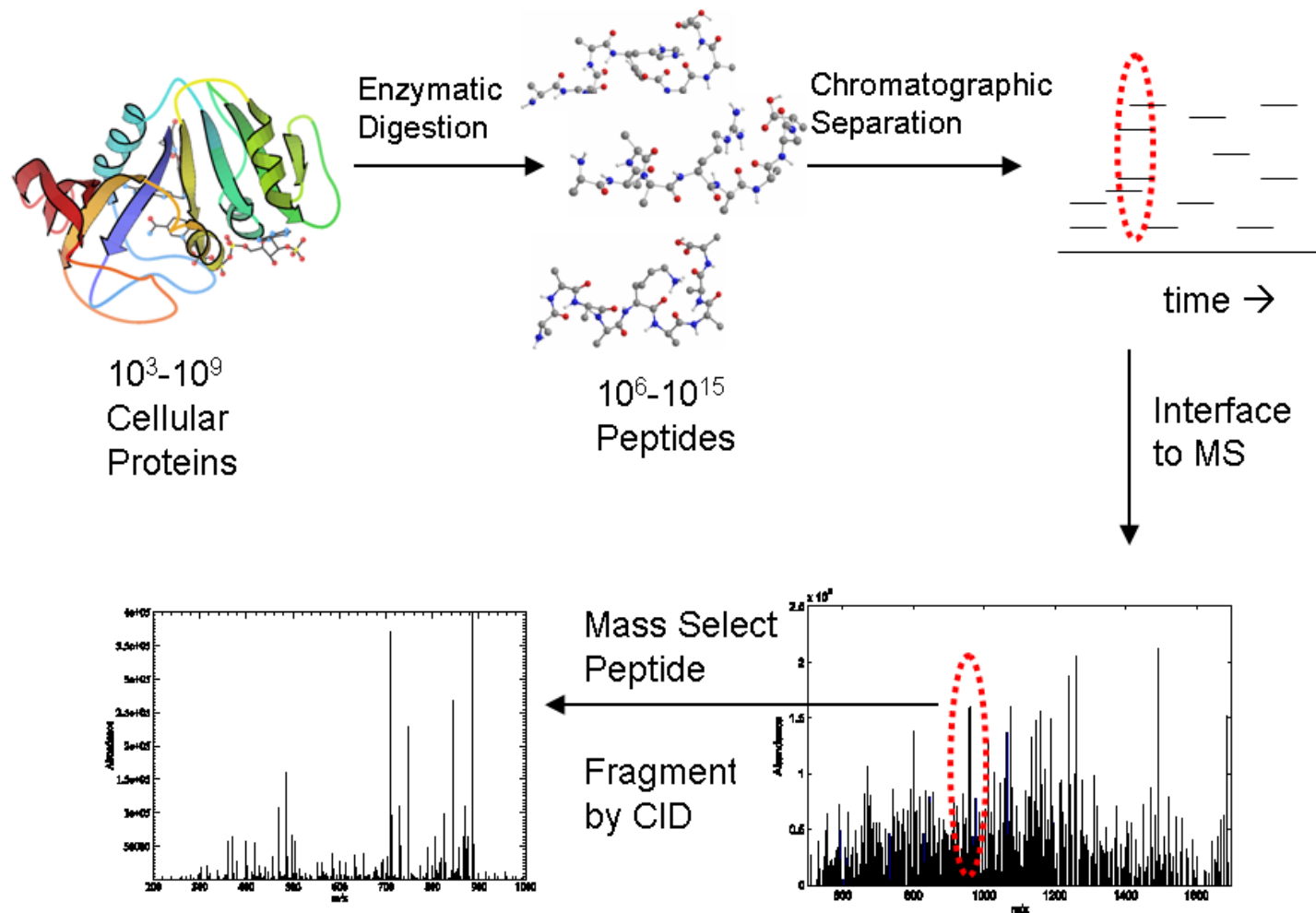
# Demand for parallel BLAST

Computing time needed to perform an all vs. all calculation grows exponentially even though compute power increases with time. Scaling demand is calculated as the number of processors required to perform an all vs. all BLAST run within 24 hours at the expected memory bandwidth capacity available at the time of the run. ScalaBLAST scales to thousands of processors, but increased scaling demand will require running on tens of thousands of processors within 3 years. Callouts indicate anticipated database size over time.
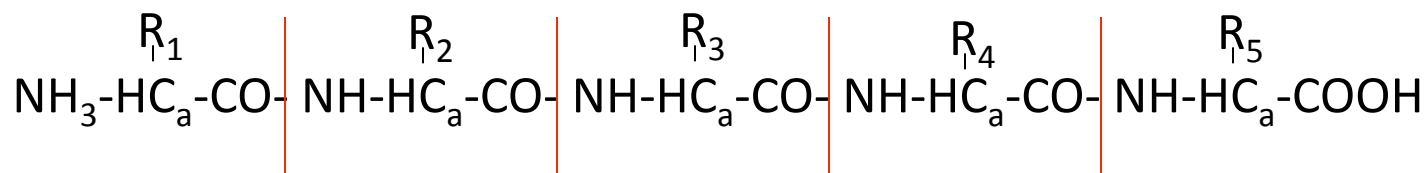
Keeping pace with sequence data

# Mass Spectrometry-based Proteomics



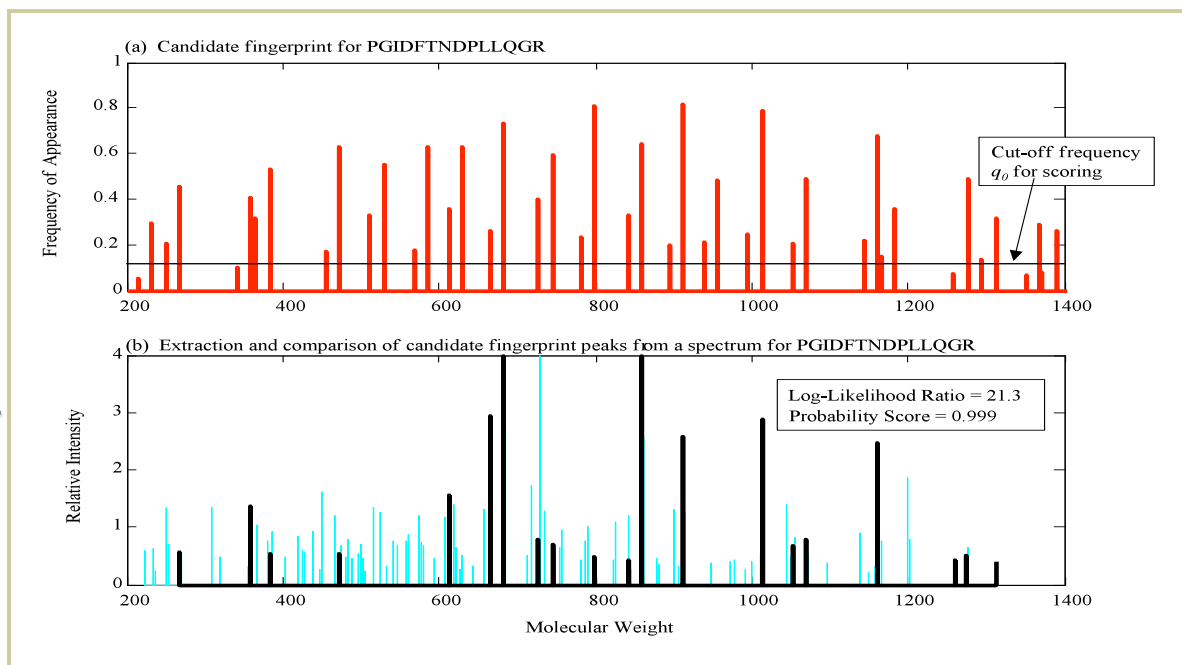Terabytes of MS-proteomics data at EMSL

# Comparing Models to Data

## Generic model spectra don't reflect the diversity of the data

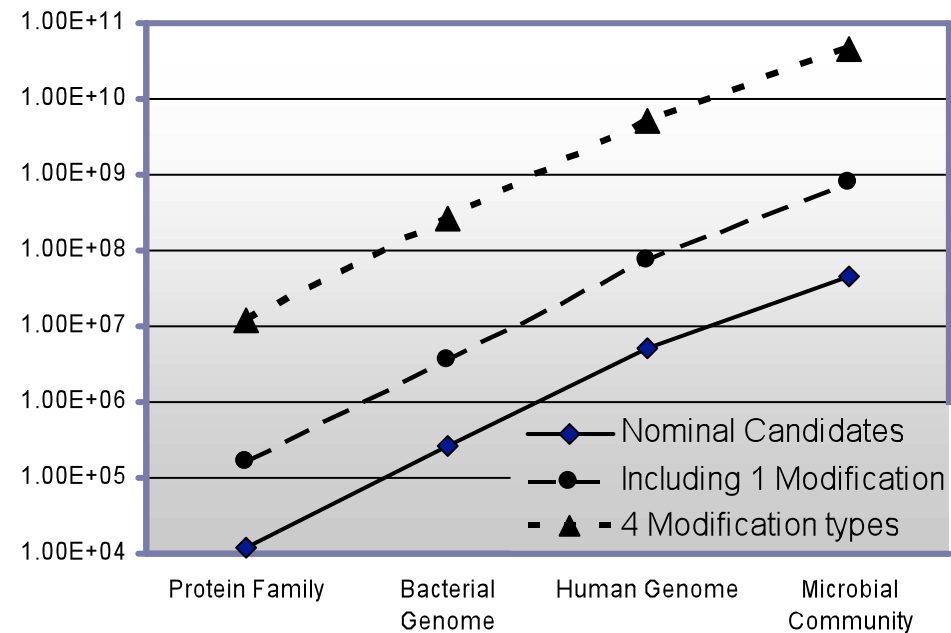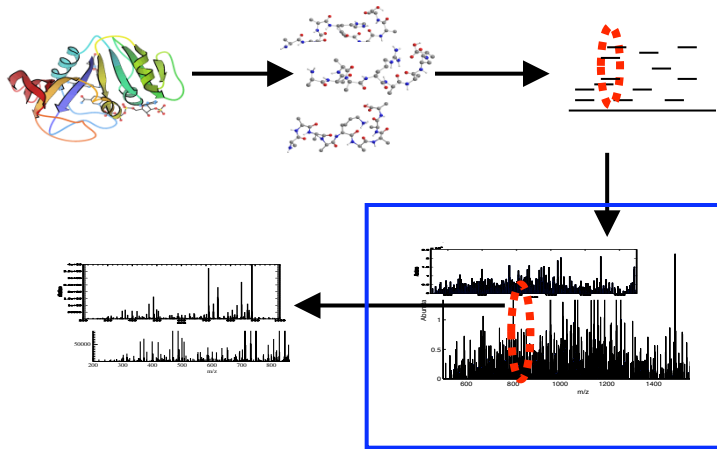$$R_1 \qquad R_2 \qquad R_3 \qquad R_4 \qquad R_5$$

$$NH_3\text{-}HC_a\text{-}CO\text{-} \;|\; NH\text{-}HC_a\text{-}CO\text{-} \;|\; NH\text{-}HC_a\text{-}CO\text{-} \;|\; NH\text{-}HC_a\text{-}CO\text{-} \;|\; NH\text{-}HC_a\text{-}COOH$$

**Model spectrum**

**Actual spectrum- Highly variable**



(a) Candidate fingerprint for PGIDFTNDPLLQGR

Frequency of Appearance

Cut-off frequency $q_0$ for scoring

(b) Extraction and comparison of candidate fingerprint peaks from a spectrum for PGIDFTNDPLLQGR

Relative Intensity

Log-Likelihood Ratio = 21.3
Probability Score = 0.999

Molecular Weight

**Models:**
- Statistical avg.
- Physical
- Experimental

# Peptide Candidates Per Spectrum



1. Not all peptides are candidate matches for each spectrum

2. Mass & chg selection

1 out of $10^5$-$10^{11}$ must be selected as the correct peptide